

# CAUSAL INFORMATION PRIORITIZATION FOR EFFICIENT REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Current Reinforcement Learning (RL) methods often suffer from sample-inefficiency, resulting from blind exploration strategies that neglect causal relationships among states, actions, and rewards. Although recent causal approaches aim to address this problem, they lack grounded modeling of reward-guided causal understanding of states and actions for goal-orientation, thus impairing learning efficiency. To tackle this issue, we propose a novel method named Causal Information Prioritization (**CIP**) that improves sample efficiency by leveraging factored MDPs to infer causal relationships between different dimensions of states and actions with respect to rewards, enabling the prioritization of causal information. Specifically, **CIP** identifies and leverages causal relationships between states and rewards to execute counterfactual data augmentation to prioritize high-impact state features under the causal understanding of the environments. Moreover, **CIP** integrates a causality-aware empowerment learning objective, which significantly enhances the agent’s execution of reward-guided actions for more efficient exploration in complex environments. To fully assess the effectiveness of **CIP**, we conduct extensive experiments across 39 tasks in 5 diverse continuous control environments, encompassing both locomotion and manipulation skills learning with pixel-based and sparse reward settings. Experimental results demonstrate that **CIP** consistently outperforms existing RL methods across a wide range of scenarios<sup>1</sup>.

## 1 INTRODUCTION

Reinforcement Learning (RL) has emerged as a powerful paradigm for training intelligent decision-making agents to learn optimal behaviors by interacting with their environments, receiving reward feedback, and iteratively optimizing their decision-making policies (Haarnoja et al., 2018; Ze et al., 2024; Sutton, 2018; Silver et al., 2017). Despite its notable successes, most RL approaches are faced with the sample-inefficiency problem, which means they typically necessitate an enormous number of interactions with the environment to learn policies, which can be impractical or costly in real-world scenarios (Savva et al., 2019; Kroemer et al., 2021). Inefficient policy learning often results from blind exploration strategies that neglect causal relationships, leading to spurious correlations and suboptimal solutions with high exploration costs (Zeng et al., 2023; Liu et al., 2024).

Causal reasoning captures essential information by analyzing causal relationships between different factors, filtering out irrelevant information, and avoiding interference from spurious correlations (Wang et al., 2022; Pitis et al., 2022; Zhang et al., 2024; Huang et al., 2022b). These approaches build internal causal structural models, enabling agents to strategically focus their exploration on the most pertinent aspects of the environment. They significantly reduce the number of samples required and demonstrate remarkable performance in single-task learning, generalization, and counterfactual reasoning (Richens & Everitt, 2024; Urpí et al., 2024; Deng et al., 2023; Huang et al., 2022a; Feng & Magliacane, 2023). However, most of these works overlook the reward-relevant causal relationships among different factors, or only partially consider the causal connections between states, actions, and rewards (Liu et al., 2024; Ji et al., 2024a), thus hindering efficient exploration.

In this work, we aim to identify and exploit task-specific causal relationships between states, actions, and rewards, enabling agents to discern relevant states and select actions that maximize rewards,

<sup>1</sup>The anonymous project page is <https://sites.google.com/view/rl-cip/>.

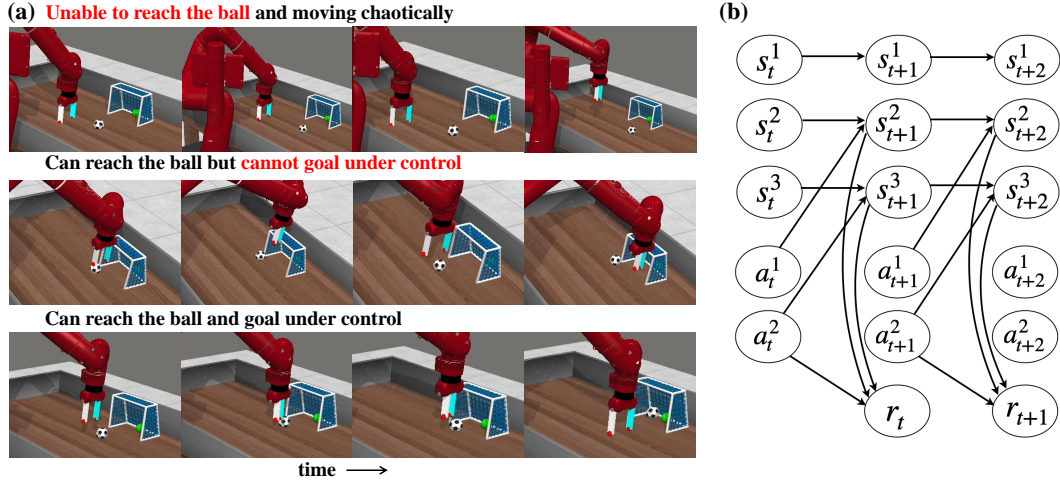


Figure 1: (a). An example of a robot manipulation soccer task with three trajectories, where the objective is to move the ball into the goal. (b). Underlying causal structure of this example in a factored MDP. Different nodes represent different dimensional states and actions.

ultimately facilitating precise and goal-oriented behaviors. Here we provide a motivating example in Figure 1, showing three trajectories for executing a manipulation soccer task, along with the underlying causal structure in a factored Markov Decision Process (MDP) (Kearns & Koller, 1999). In the first trajectory (row 1), when the agent fails to distinguish states with more intricate causal relationships of the task, the robotic arm exhibits chaotic moving and receives no rewards. The second trajectory (row 2) shows that even without chaotic movements, uncontrollable actions unrelated to the reward lead to an inability to guide the ball towards the goal. Only by filtering out irrelevant state features and executing more controllable actions can we guarantee that the ball is kicked into the goal like row 3. Quantifying the contribution of different factors to the reward can effectively help analyzing important causal relationships.

To address the limitation of sample-inefficiency and leverage the potential of causal reasoning, we propose a novel approach named Causal Information Prioritization (CIP) for efficient RL, improving learning efficiency from the perspective of rewards. Building upon the factored MDPs, CIP infers causal relationships between states, actions, and rewards across different dimensions, respectively. CIP employs counterfactual data augmentation based on the causality between states and rewards to generate transitions, prioritizing critical state transitions. Furthermore, CIP leverages the causality between actions and rewards to reweight actions, while utilizing empowerment to maximize mutual information between causally informed actions and future states, thereby enabling better control.

Specifically, CIP leverages collected data to construct a reward-guided structural model that explicitly reasons about state-reward causal influences, enabling the swapping of causally independent state features across observed trajectories to generate synthetic transitions without additional environment interactions. By swapping independent state features across different transitions (i.e., irrelevant state dimensions of chaotic movements in the soccer task), CIP accentuates causally dependent state information (i.e., relevant states to reach the ball), facilitating focused learning of critical state transitions. Subsequently, CIP constructs another structural model that incorporates actions and rewards to reweight actions of dimensions. To enhance the exploration efficiency, CIP integrates a causality-aware empowerment term, quantifying the agent’s capacity to exert controlled influence over its environment through the mutual information. This empowerment term, combined with causally weighted actions, is integrated into the learning objective, prioritizing actions with high causal influence. The synthesis of causal reasoning and action empowerment enables agents to focus on behaviors that are causally relevant to the task, leading to more efficient and effective policy learning. The main contributions of this work can be summarized as follows.

- To address limitations of blind exploration and sample-inefficiency, we introduce CIP, a novel efficient RL framework that prioritizes causal information through the lens of reward. CIP bridges the gap between causal reasoning and empowerment to facilitate efficient exploration.

- **CIP** constructs reward-guided structural models to uncover causal relationships between states, actions, and rewards across dimensions. By leveraging state-reward causality, it performs counterfactual data augmentation, eliminating the need for additional environment interactions, and enabling learning on critical state transitions. Exploiting action-reward causality, it reweights actions to enhance exploration efficiency through empowerment. By prioritizing causal information, **CIP** enables agents to focus on behaviors that have causally significant effects on their tasks.
- To validate the effectiveness of **CIP**, we conduct extensive experiments in 39 tasks across 5 diverse continuous control environments, including manipulation and locomotion. These comprehensive evaluations demonstrate the effectiveness of **CIP** in pixel-based and sparse reward settings, underscoring its versatility and reliability.

## 2 RELATED WORK

### 2.1 CAUSAL RL

The application of causal reasoning in RL has shown significant potential to improve sample efficiency and generalization by effectively excluding irrelevant environmental factors through causal analysis (Huang et al., 2022a; Feng & Magliacane, 2023; Mutti et al., 2023). Wang (Wang et al., 2021) introduces a novel regularization-based method for causal dynamics learning, which explicitly identifies causal dependencies by regulating the number of variables used to predict each state variable. CDL (Wang et al., 2022) takes an innovative approach by using conditional mutual information to compute causal relationships between different dimensions of states and actions, explicitly eliminating irrelevant components. IFactor (Liu et al., 2024) is a general framework to model four distinct categories of latent state variables, capturing various aspects of information. ACE (Ji et al., 2024a), an off-policy actor-critic method, integrates causality-aware entropy regularization. Although ACE makes significant strides in uncovering causal relationships between actions and rewards, it notably overlooks the role of states in these causal structures. Table 2 provides a categorization of various causal RL methods, highlighting their focus on different reward-guided causal relationships. Existing approaches do not fully account for the causal relationships between both states and actions with rewards, and their application to more general tasks remains limited. Our goal is to comprehensively explore these causal relationships from a reward-guided perspective to enhance sample efficiency across a broader range of tasks.

### 2.2 EMPOWERMENT IN RL

Empowerment, an information theory-based concept of intrinsic motivation, has emerged as a powerful paradigm for enhancing an agent’s environmental controllability (Mohamed & Jimenez Rezende, 2015; Klyubin et al., 2005; Cao et al., 2024). This framework conceptualizes actions and future states as information transmission channels, offering a novel perspective on agent-environment interactions. In RL, empowerment has been applied to uncover more controllable associations between states and actions, as well as to develop robust skill (Salge et al., 2014; Bharadhwaj et al., 2022; Choi et al., 2021; Eysenbach et al., 2018; Leibfried et al., 2019; Seitzer et al., 2021). Empowerment, expressed as maximizing mutual information  $\max_{\pi} I$ , serves as a learning objective in various RL frameworks, providing intrinsic motivation for exploration and potentially yielding more efficient and generalizable policies. Our approach extends empowerment in RL by examining the influence of state, actions, and rewards through a causal lens, integrating causal understanding with empowerment to enhance exploration strategy and learning efficiency.

## 3 PRELIMINARIES

### 3.1 MARKOV DECISION PROCESS

In RL, the agent-environment interaction is formalized as an MDP. The standard MDP is defined by the tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mu_0, r, \gamma \rangle$ , where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  represents the action space,  $\mathcal{P}(s'|s, a)$  is the transition dynamics,  $r(s, a)$  is the reward function, and  $\mu_0$  is the distribution of the initial state  $s_0$ . The discount factor  $\gamma \in [0, 1]$  is also included. The objective of RL is to learn a policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  that maximizes the expected discounted cumulative reward  $\eta_{\mathcal{M}}(\pi) := \mathbb{E}_{s_0 \sim \mu_0, s_t \sim \mathcal{P}, a_t \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ .

### 3.2 STRUCTURAL CAUSAL MODEL

A Structural Causal Model (SCM) (Pearl, 2009) is defined by a distribution over random variables, defined as  $\mathcal{V} = \{s_t^1, \dots, s_t^d, a_t^1, \dots, a_t^n, r_t, s_{t+1}^1, \dots, s_{t+1}^d\}$  and a Directed Acyclic Graph (DAG)  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with a conditional distribution  $\mathcal{P}(v_i | \text{PA}(v_i))$  for node  $v_i \in \mathcal{V}$ . Then the distribution can be specified as:

$$p(v_1, \dots, v_{|\mathcal{V}|}) = \prod_{i=1}^{|\mathcal{V}|} p(v_i | \text{PA}(v_i)), \quad (1)$$

where  $\text{PA}(v_i)$  is the set of parents of the node  $v_i$  in the graph  $\mathcal{G}$ .

**Causal Structures in MDP** We use a factored MDP (Kearns & Koller, 1999; Guestrin et al., 2003; 2001) to model the MDP and the underlying causal structures between states, actions, and rewards. In the factored MDP, nodes represent system variables (rewards and different dimensions of the states and actions), while the edges denote their relationships within the MDP. We employ causal discovery methods to learn the structures of  $\mathcal{G}$ .

We can identify the graph structure in  $\mathcal{G}$ , which can be represented as the adjacency matrix  $M$ . To integrate such relationships in MDP, we explicitly encode the causal mask over variables into the reward function. Hence, the reward function in MDP with the causal structure is defined as follows:

$$r_t = R(M^{s \rightarrow r} \odot s_t, M^{a \rightarrow r} \odot a_t, \epsilon_{r,t}) \quad (2)$$

where  $\odot$  denotes the element-wise product.  $M^{s \rightarrow r} \in \mathbb{R}^{|s| \times 1}$  and  $M^{a \rightarrow r} \in \mathbb{R}^{|a| \times 1}$  are the adjacency matrices indicating the influence of current states and actions on the reward, respectively, and  $\epsilon_{r,t}$  represents i.i.d. Gaussian noise. Under the Markov condition and faithfulness assumption (Pearl, 2009; Spirtes et al., 2001), the structural vectors are identifiable. The detailed assumptions and propositions can be found in Appendix B. In this work, our objective is to discover and leverage these two causal matrices to prioritize causal information for efficient RL.

### 3.3 EMPOWERMENT IN RL

Empowerment quantifies an agent’s capacity to influence its environment and perceive the consequences of its actions (Klyubin et al., 2005; Bharadhwaj et al., 2022; Jung et al., 2011). In our framework, the empowerment is defined as the mutual information between the agent state  $s_{t+1}$  and action  $a_t$ , conditioned on the present state  $s_t$  and causal mask  $M$ , as shown follows:

$$\mathcal{E} := \max_{p(a_t)} \mathcal{I}(a_t; s_{t+1} \mid s_t, M), \quad (3)$$

where  $\mathcal{E}$  denotes the channel capacity from action to state, and  $p(a_t)$  is the distribution of actions. Unlike (Cao et al., 2024), which focuses on action-to-state empowerment effects, we leverage causal understanding and more accurate entropy calculation to analyze state-to-action influences, facilitating the development of more controllable behavioral policies.

## 4 CIP

In this section, we introduce the proposed framework **CIP**, which implements causal information prioritization based on the causal relationships between states, actions, and rewards (as shown in Figure 2). First, we train a structural model based on the causal discovery method, DirectLiNGAM (Shimizu et al., 2011) using collected trajectories to obtain a causal matrix  $M^{s \rightarrow r}$ . Utilizing this matrix, **CIP** executes the swapping of causally independent state features, generating synthetic transitions (Section 4.1). This process of swapping independent state information accentuates causally dependent state information, enabling focused learning on critical state transitions. Subsequently, **CIP** constructs another structural model to get a weight matrix  $M^{a \rightarrow r}$  that incorporates actions and rewards to reweight actions (Section 4.2). Furthermore, **CIP** integrates a causality-aware empowerment term  $\mathcal{E}_{\pi_c}(s)$  combined with causally weighted actions into the learning objective to promote efficient exploration. This integration encourages the agent’s policy  $\pi_c$  to prioritize actions with high causal influence, thereby enhancing its goal-achievement capabilities.



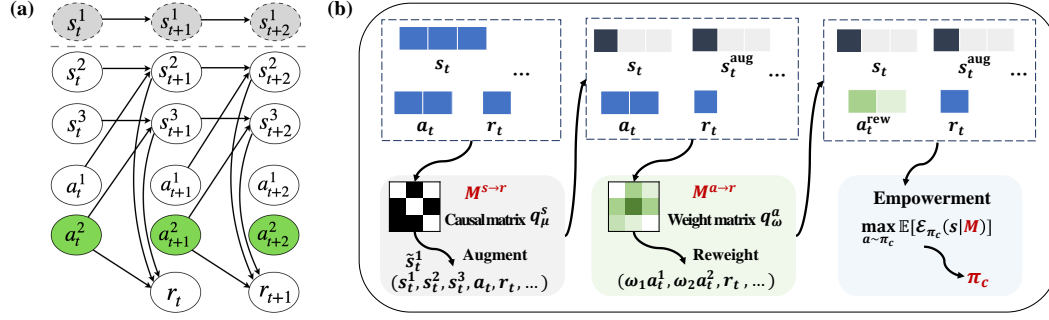


Figure 2: (a) Underlying causal structure of **CIP**. (b) The whole learning process of **CIP** includes counterfactual data augmentation, causal action reweight and causal action empowerment.

#### 4.1 COUNTERFACTUAL DATA AUGMENTATION

To discover the causal relationships between states and rewards, we initially collect trajectories to train a structural model by the DirectLiNGAM method, denoted as  $q_\mu^s$ , to obtain the causal matrix  $M^{s \rightarrow r}$ . Subsequently, we infer the local factorization, which is utilized to generate counterfactual transitions. For each state  $s$  in the trajectories, we compute the uncontrollable set, defined as the set of variables in  $s$  for which the agent has no causal influence on rewards:

$$\mathcal{U}_s = \{s^i \mid M^{s \rightarrow r} \cdot (s^i, r_t) < \theta; i \in [1, N]\}, \quad (4)$$

where  $\theta$  is a fixed threshold and  $N$  is the dimension of the state space. The set  $\mathcal{U}_s$  encompasses all dimensional state variables for which the causal relationship  $s_t^i \rightarrow r_t$  does not exist in the causal matrix of states and rewards. Utilizing the learned causal matrix  $M^{s \rightarrow r}$ , we partition all state variables in the factored MDP into controllable and uncontrollable sets. These uncontrollable sets are then leveraged for counterfactual data augmentation, thereby prioritizing the causally-informed state information to improve learning efficiency.

To generate counterfactual samples, we perform a swap of variables that fall under the uncontrollable category (i.e., in set  $\mathcal{U}_s$ ) sampled from the collect trajectories. Specifically, given two transitions  $(s_t, a_t, s_{t+1}, r_t)$  and  $(\hat{s}_j, \hat{a}_j, \hat{s}_{j+1}, \hat{r}_j)$  sampled from trajectories, which share at least one uncontrollable sub-graph structure (i.e.,  $\mathcal{U}_s \cap \mathcal{U}_{\hat{s}} \neq \emptyset$ ), we construct a counterfactual transition  $(\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1}, \tilde{r}_t)$  by swapping the irrelevant state variables  $(s_t^i, s_{t+1}^i)$  with  $(\hat{s}_j^i, \hat{s}_{j+1}^i)$  for each  $i \in \mathcal{U}_s \cap \mathcal{U}_{\hat{s}}$ . The augmented transitions will be added to the training data for causal reasoning during subsequent action empowerment, thus eliminating the need for additional environment interactions to prioritize causal information. Furthermore, we also consider directly using controllable state sets combined with causal action empowerment to replace counterfactual data augmentation for policy learning. The comparative experimental results validating this approach are presented in Appendix D.3.1.

#### 4.2 CAUSAL ACTION PRIORITIZATION THROUGH EMPOWERMENT

**Causal action reweight** Having analyzed the causal relationships between states and rewards to achieve efficient data augmentation, in this section, we further discover the causal relationships between actions and rewards to prioritize causally-informed decision-making behaviors. **CIP** constructs a reward-guided structural model, incorporating states (including augmented states), actions, and rewards. This model forms the foundation for action prioritization in policy learning, enabling action reweighting based on causality. Leveraging this structural model to delineate relationships between policy decisions and rewards, we evaluate the causal impact of different actions on reward outcomes. In this way, the agent focuses on pivotal actions with demonstrable causal links to desired reward outcomes, potentially accelerating learning and optimizing performance in complex environments.

Specifically, in **CIP**, we employ DirectLiNGAM method to train a causal structural model  $q_\omega^a$ , which yields a weight matrix  $M^{a \rightarrow r}$ , delineating the relationships between actions and rewards, conditioned on the states. For a given set of actions  $(a_t^1, a_t^2, a_t^3, \dots)$ , we utilize the weight matrix  $M^{a \rightarrow r}$  to reweight them as  $(\omega_1 a_t^1, \omega_2 a_t^2, \omega_3 a_t^3, \dots)$ , where  $\omega$  represents the causal weights derived from the matrix  $M^{a \rightarrow r}$ . By leveraging this causal structure, we can prioritize the most pivotal actions, potentially leading to more efficient policy exploration and targeted policy improvements.

**Causal action empowerment** Based on the learned causal structure, we propose the causal action empowerment to incorporate the reweighted actions into the learning objective for efficient exploration in a controllable manner. To this end, we design a causality-aware empowerment term  $\mathcal{E}_{\pi_c}(s)$  for policy optimization. We maximize the empowerment gain of the policy  $\pi_c$ , where  $\pi_c$  incorporates the learned causal structure. This approach allows us to quantify and maximize the empowerment that can be achieved by explicitly considering causal relationships, thereby bridging the gap between causal reasoning and empowerment.

We denote the empowerment of the causal policy as  $\mathcal{E}_{\pi_c}(s) = \max_a \mathcal{I}(a_t; s_{t+1} | s_t; M)$ . We then formulate the following objective empowerment function:

$$\begin{aligned} \mathcal{E}_{\pi_c}(s) &= \max_a \mathcal{I}(a_t; s_{t+1} | s_t; M) \\ &= \max_{a_t \sim \pi_c(\cdot | s)} \mathcal{H}(\pi_c(a_t | s_t)) - \mathcal{H}(\pi_c(a_t | s_t; s_{t+1})), \end{aligned} \quad (5)$$

where  $\pi_c$  is the policy under the causal weighted matrix  $M^{a \rightarrow r}$ . The first entropy term  $\mathcal{H}(\pi_c(a_t | s_t))$  promotes action diversity within the constraints of the causal structure. It encourages the agent to explore a wide range of actions that are causally informed, while the second entropy term  $-\mathcal{H}(\pi_c(a_t | s_t; s_{t+1}))$  enhances the action predictability in state transitions. It encourages the selection of actions that lead to predictable outcomes, given the current and subsequent states, thereby promoting controlled and goal-oriented behaviors. We train an inverse dynamics model to represent the policy  $\pi_c(\cdot | s_t; s_{t+1})$ . The detailed derivation proceeds as follows:

$$\mathcal{H}(\pi_c(\cdot | s_t)) = -\mathbb{E}_{a_t \in \mathcal{A}} \left[ \sum_{i=1}^{d_A} M^{a^i \rightarrow r} \odot \pi_c(a_t^i | s_t) \log \pi(a_t^i | s_t) \right], \quad (6)$$

and

$$\mathcal{H}(\pi_c(\cdot | s_t; s_{t+1})) = -\mathbb{E}_{a_t \in \mathcal{A}} \left[ \sum_{i=1}^{d_A} M^{a^i \rightarrow r} \odot \pi_c(a_t^i | s_t; s_{t+1}) \log \pi(a_t^i | s_t; s_{t+1}) \right], \quad (7)$$

where  $d_A$  is the dimension of the action space. Hence, the learning objective of the causal action empowerment can be defined as follows:

$$\begin{aligned} \mathcal{E}_{\pi_c}(s) &= \max_{a_t \sim \pi_c(\cdot | s)} \mathcal{H}(\pi_c(a_t | s_t)) - \mathcal{H}(\pi_c(a_t | s_t; s_{t+1})) \\ &= \max_{a_t \sim \pi_c(\cdot | s)} \mathbb{E}_{\pi_c(a_t | s_t) p_{\pi_c}(a_t | s_t, s_{t+1})} [\log \mathcal{P}_{\phi_c}(a_t | s_t, s_{t+1}; M) - \log \mathcal{P}_{\pi_c}(a_t | s_t; M)], \end{aligned} \quad (8)$$

where  $\mathcal{P}_{\pi_c}(a_t | s_t; M)$  is the action distribution given current state of policy  $\pi_c$  with the causal structure, which can be denoted as  $\pi_c(a_t | s_t)$ .  $\mathcal{P}_{\phi_c}(a_t | s_{t+1}, s_t; M)$  represents an inverse dynamics model trained on the collected transitions of state variables. Hence, we update the target policy  $\pi_c$  by maximizing the empowerment objective function derived in Eq. 8.

Adhering to the maximum entropy paradigm (Haarnoja et al., 2018), we calculate  $\mathcal{E}_{\pi_c}(s)$  for maximization instead of standard entropy, thus prioritizing exploration of pivotal actions that are more likely to have significant causal effects on the reward. This targeted exploration strategy has the potential to accelerate learning by focusing on the most influential actions in current controllable states. Based on the causality-aware empowerment, the Q-value for policy  $\pi_c$  could be computed iteratively by applying a modified Bellman operator  $\mathcal{T}_c^\pi$  with  $\mathcal{E}_{\pi_c}(s)$  term as stated below:

$$\begin{aligned} \mathcal{T}_c^\pi Q(s_t, a_t) &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_t \sim \pi_c} [Q(s_{t+1}, a_{t+1}) + \alpha \mathcal{E}_{\phi_c}(s)]] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_t \sim \pi_c} [Q(s_{t+1}, a_{t+1}) + \alpha (\mathcal{H}(\pi_c(a_t | s_t)) - \mathcal{H}(\pi_c(a_t | s_t; s_{t+1})))]]. \end{aligned} \quad (9)$$

Hence, we integrate the causality-aware empowerment term into the policy optimization objective function,  $\hat{\eta}_{\mathcal{M}}(\pi_c) = \mathbb{E}_{s_0 \sim \mu_0, s_t \sim \mathcal{P}, a_t \sim \pi_c} [\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \alpha \mathcal{E}_{\pi_c}(s))]$ .

In summary, **CIP** harnesses empowerment to integrate the causal understanding into decision-making. By maximizing the empowerment gain of the causally-informed policy, we guide the agent to prioritize actions that align with the environment’s underlying causal relationships. This approach enhances the agent’s exploration efficiency, focusing on actions with meaningful causal impacts and correlated with desired outcomes. Algorithm 1 illustrates the complete **CIP** pipeline.

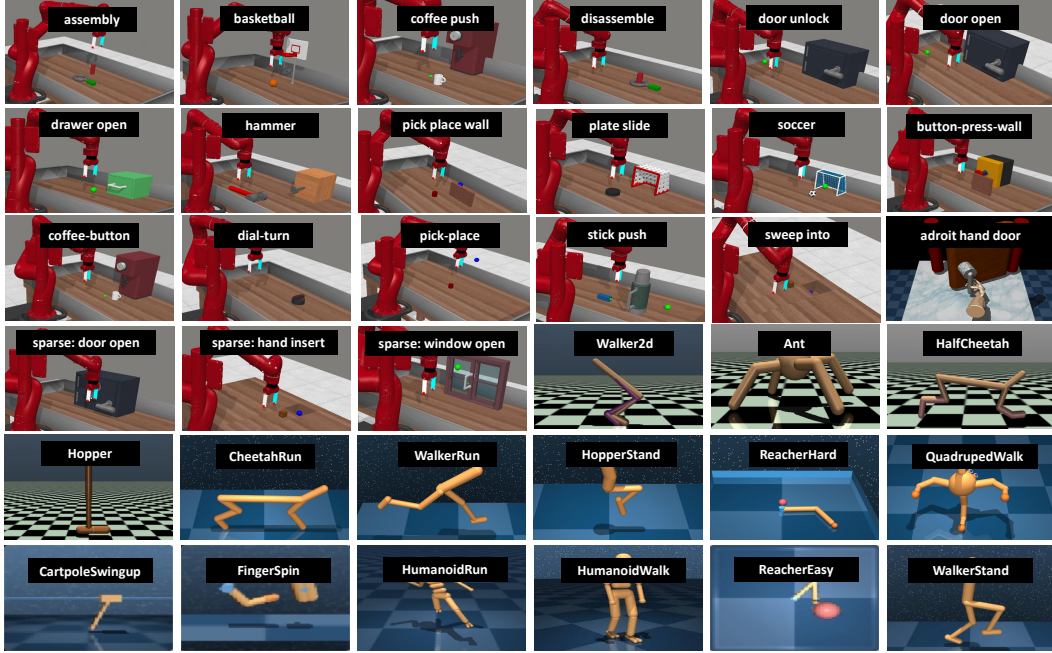


Figure 3: The 36 experimental tasks in 5 continuous control environments

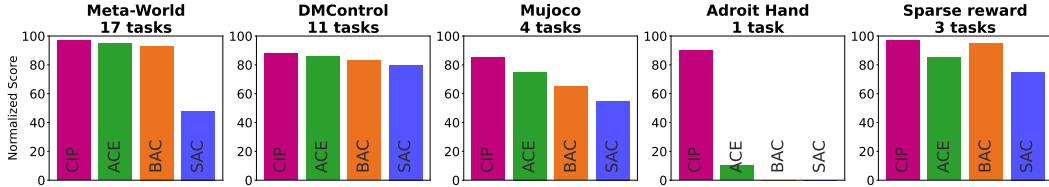


Figure 4: Experimental results with normalized score across all 36 tasks in 5 environments.

## 5 EXPERIMENTS

Our experiments aim to address the following questions: (i) How does the performance of **CIP** compare to other RL approaches in diverse continuous control tasks, including manipulation and locomotion with sparse rewards, high-dimensional action spaces, and pixel-based challenges? (ii) Can **CIP**, through data augmentation and empowerment, improve sample efficiency and learn reliable policies? (iii) What are the effects of the components and hyperparameters in **CIP**?

### 5.1 EXPERIMENTAL SETUP

**Environments.** We evaluate **CIP** on 5 continuous control environments, including MuJoCo (Todorov et al., 2012), DMControl (Tassa et al., 2018), Meta-World (Yu et al., 2020), Adroit Hand (Rajeswaran et al., 2018), and sparse reward setting environments in Meta-World. This comprehensive evaluation encompasses 36 tasks, spanning both locomotion and manipulation skill learning, as illustrated in Figure 3. We also conduct experiments in 3 pixel-based tasks of the DMControl environment as shown in Figure 17. Our experimental tasks incorporate a wide range of challenges, including high-dimensional state and action spaces, sparse reward settings, pixel-based scenarios, and locomotion. For extensive experimental settings, please refer to Appendix D.1.

**Baselines.** We compare **CIP** with three popular RL baselines across all 36 tasks and against IFactor (Liu et al., 2024) in 3 pixel-based tasks: (1) SAC (Haarnoja et al., 2018), an off-policy actor-critic algorithm featuring maximum entropy regularization. (2) ACE (Ji et al., 2024a), a method employing causality-aware entropy regularization. (3) BAC (Ji et al., 2024b), a method that balances sufficient exploitation of past successes with exploration optimism. (4) IFactor (Liu et al., 2024), a causal framework modeling four distinct categories of latent state variables for pixel-based tasks. To ensure robustness and statistical significance, we conduct each experiment using 4 random seeds.

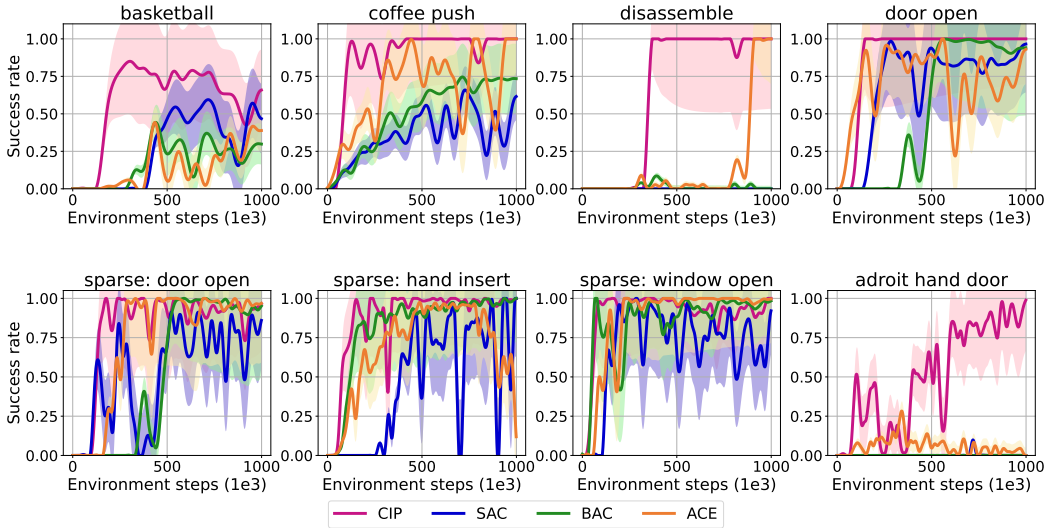


Figure 5: Experimental results of 8 manipulation skill learning tasks in Meta-World and adroit hand environments including sparse reward settings. For all tasks results, please refer to Appendix D.2.

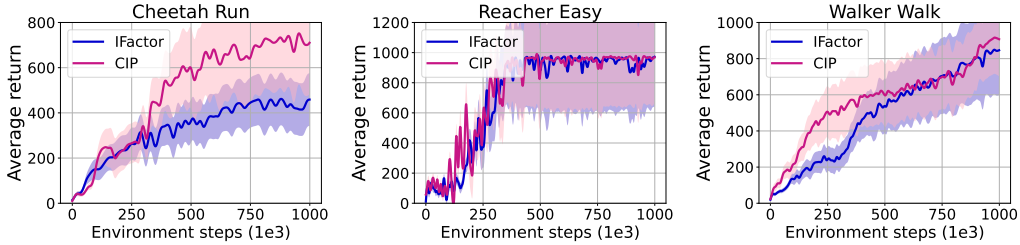


Figure 6: Experimental results of 3 pixel-based tasks in DMControl environment.

## 5.2 MAIN RESULTS

Figure 4 presents the normalized scores of **CIP** compare to other methods across 36 tasks in 5 environments. In 17 Meta-World robot-arm tasks, **CIP** achieves a near-perfect score of 100, showcasing its exceptional performance in manipulation tasks. For locomotion tasks in DMControl and MuJoCo, **CIP** consistently attains scores exceeding 80, indicating robust performance across diverse locomotion challenges. Notably, **CIP** exhibits significant performance improvements in challenging scenarios, such as adroit hand manipulation and 3 tasks with the sparse reward setting. These results underscore the effectiveness in tackling complex, high-dimensional control problems. In next sections, we present a comprehensive analysis of **CIP**'s performance across diverse tasks.

**Robot-arm manipulation.** Figure 5 presents the success rates across 7 Meta-World robot-arm manipulation tasks including sparse reward settings. **CIP** consistently outperforms all other methods across these tasks, demonstrating both faster policy learning and enhanced stability. In challenging tasks, such as disassemble, **CIP** achieves an impressive 100% success rate. The effectiveness of **CIP** can be attributed to focus on causally relevant information within the state and action spaces. In sparse reward settings, the efficient extraction of causal state information and the prioritization of controllable actions enable effective task completion. By systematically eliminating noise from non-causal factors, **CIP** allows the agent to construct a more controllable and efficient policy, directly addressing the core challenges in the given task domain.

**High-dimensional Adroit hand manipulation.** To rigorously evaluate our method's efficacy in high-dimensional tasks, we conduct comparative experiments in the Adroit Hand environment of door open task. This challenging setup involves controlling a robotic hand with up to 28 actuated degrees of freedom ( $\mathcal{A} \in \mathbb{R}^{28}$ ). Figure 5 illustrates the success rates achieved across all methods. Notably, while the three other comparative methods fail to demonstrate significant progress on this challenging task, **CIP** achieves a near 100% success rate after 700k environment steps. **CIP** demonstrates its capacity



Table 1: The comparative experimental results of average return in 8 locomotion tasks. We bold the best scores, and underline suboptimal results,  $\pm$  is the standard deviation, w/o represents without.

Method	Ant	HalfCheetah	Hopper	Walker2d	Cheetah Run	Hopper Stand	Quadruped Walk	Reacher Hard
<b>CIP</b>	6418 $\pm$ 81	<b>12594<math>\pm</math>210</b>	<b>2846<math>\pm</math>882</b>	<b>5624<math>\pm</math>91</b>	<b>893<math>\pm</math>12</b>	<b>936<math>\pm</math>17</b>	948 $\pm$ 54	<b>991<math>\pm</math>11</b>
w/o Aug	6231 $\pm$ 81	<u>12225<math>\pm</math>102</u>	2308 $\pm$ 785	5294 $\pm$ 41	<u>885<math>\pm</math>13</u>	931 $\pm$ 22	945 $\pm$ 35	<u>989<math>\pm</math>13</u>
w/o Emp	6295 $\pm$ 210	10986 $\pm$ 572	2270 $\pm$ 904	<u>5547<math>\pm</math>91</u>	876 $\pm$ 21	785 $\pm$ 114	924 $\pm$ 23	971 $\pm$ 13
<b>SAC</b>	6062 $\pm$ 105	10888 $\pm$ 240	2266 $\pm$ 981	5251 $\pm$ 106	767 $\pm$ 16	<b>936<math>\pm</math>8</b>	930 $\pm$ 19	980 $\pm$ 8
<b>BAC</b>	<b>6511<math>\pm</math>30</b>	10276 $\pm$ 34	2263 $\pm$ 1063	3316 $\pm$ 702	665 $\pm$ 6	932 $\pm$ 4	<b>962<math>\pm</math>24</b>	974 $\pm$ 16
<b>ACE</b>	5922 $\pm$ 106	9390 $\pm$ 25	<u>2312<math>\pm</math>673</u>	4922 $\pm$ 96	863 $\pm$ 23	912 $\pm$ 16	933 $\pm$ 57	973 $\pm$ 17

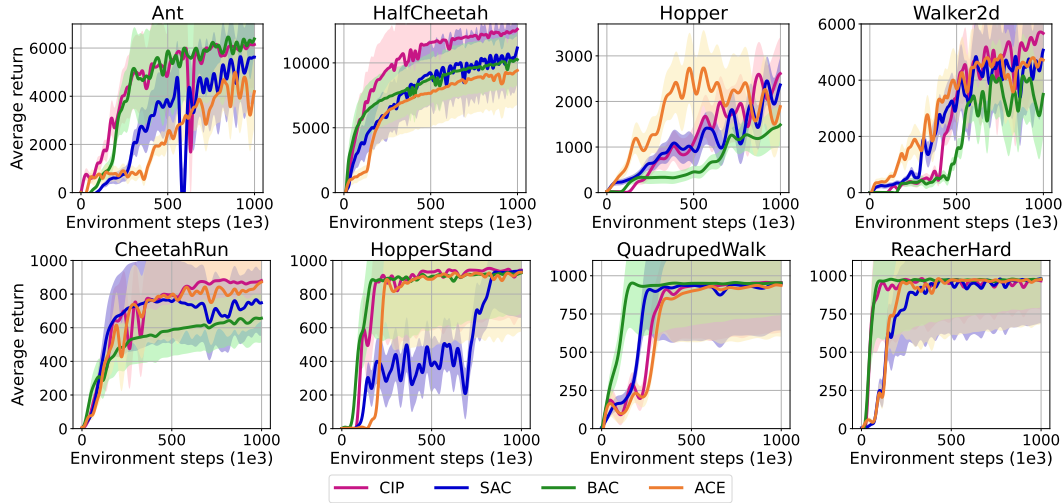


Figure 7: Experimental results with average return across 8 tasks in locomotion tasks.

to identify and focus on the most pertinent variables, even within highly complex environments. The proposed framework significantly improves the policy’s ability to exert precise control.

**Locomotion.** We further evaluate **CIP** in another important category: locomotion. The part experimental results of average return in MuJoCo and DMControl environments are presented in Table 1. Learning curves are illustrated in Figure 7. We observe that **CIP** achieves the best performance in six tasks and sub-optimal in other tasks. Moreover, compared to the traditional method SAC, **CIP** demonstrates significant performance improvements in more challenging tasks such as CheetahRun and Hopper. Compared to the causality-based method ACE, **CIP** demonstrates improvements in all tasks. Overall, in locomotion tasks, **CIP** achieves superior performance and attains high sample efficiency. Detailed performance on all benchmarks are provided in Appendix D.2.

**Pixel-based task learning** To further validate the performance in pixel-based tasks, we use 3 complex pixel-based DMControl tasks for evaluation, where video backgrounds serve as distractors. We apply the proposed counterfactual data augmentation and causal action empowerment to the IFactor method for comparative analysis. As shown in Figure 6, **CIP** surpasses IFactor in terms of average return. These results underscore **CIP**’s efficacy in pixel-based tasks and its capacity to better overcome spurious correlations arising from video backgrounds, focusing on the locomotion. For visualization trajectories in pixel-based results, please refer to Appendix D.2.4.

### 5.3 PROPERTY ANALYSIS

**Ablation study.** We conduct ablation experiments involving **CIP**, **CIP** without (w/o) counterfactual data augmentation (Aug), and **CIP** w/o Empowerment (Emp). The results in 8 locomotion tasks are shown in Table 1. And the learning curves of all tasks are depicted in Appendix D.3. The experimental results reveal that the variant without the empowerment learning objective performs poorly, underscoring the critical role of empowerment maximization in enhancing control capabilities.



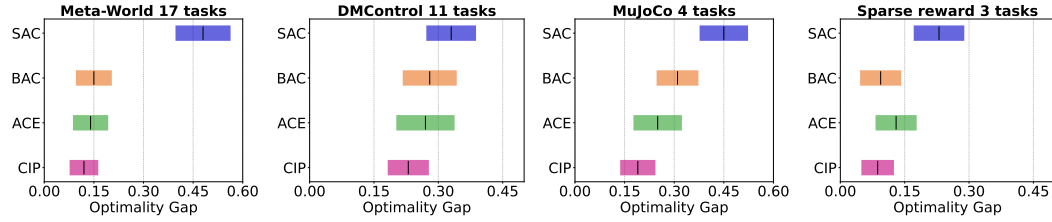


Figure 8: Experimental results of reliability evaluation by the metric Optimality Gap (lower values are better) on 4 diverse environments across 35 tasks.

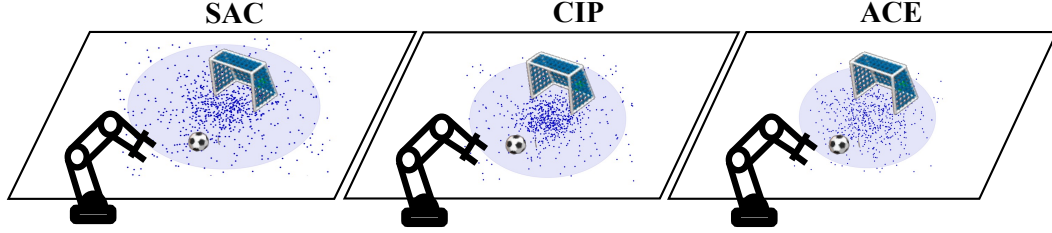


Figure 9: Visualization of the trajectories in soccer task.

Additionally, **CIP** without counterfactual data augmentation is less sample efficient than **CIP**, highlighting the importance of augmentation.

**Reliability evaluation.** We evaluate **CIP**’s reliability across 35 tasks in 4 environments, excluding the Adroit Hand door task due to **CIP**’s exceptional performance there. Figure 8 illustrates the experimental results using the Optimality Gap metric (Agarwal et al., 2021). **CIP** consistently achieves the lowest values across all tasks in four environments, with lower values indicating superior performance. This consistent excellence across diverse scenarios underscores the robustness and reliability of our proposed method. These results highlight **CIP**’s effectiveness in enhancing agent performance across a wide range of tasks, demonstrating its versatility and broad applicability.

**Visualization.** We employ trajectory visualization to comparatively validate the efficacy of our method. As depicted in Figure 9, the light-shaded regions delineate the policy exploration space, while the point clustering area indicates the area of frequent interaction. Our analysis reveals that **CIP**, leveraging counterfactual data augmentation, achieves substantially broader exploration compared to ACE and SAC. Concurrently, the causal information prioritization framework facilitates more focused execution in critical state regions. These visual findings provide robust empirical support for the effectiveness of our proposed augmentation framework.

## 6 CONCLUSION

This study introduces an efficient RL framework, designed to enhance sample efficiency. This approach begins by counterfactual data augmentation using the causality between states and rewards, effectively mitigating interference from irrelevant states without additional environmental interactions. We then develop a reward-guided structural model that leverages causal awareness to prioritize causal actions through empowerment. We conduct extensive experiments across 39 tasks spanning 5 diverse continuous control environments which demonstrate the exceptional performance of our proposed method, showcasing its robustness and adaptability across challenging scenarios.

**Limitation** In this paper, **CIP** considers combining model-free RL algorithms. To further enhance learning efficiency and generalization, it’s worthwhile to construct a causal world model capable of capturing accurate state transitions, which requires additional consideration of the causal relationships between states and actions. Addressing this issue will be the focus of our future work.

## REPRODUCIBILITY STATEMENT

We provide the core code of **CIP** in the supplementary material. The implementation details are shown in Appendix D.1.

## REFERENCES

- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.
- Homanga Bharadhwaj, Mohammad Babaeizadeh, Dumitru Erhan, and Sergey Levine. Information prioritization through empowerment in visual model-based rl. In *International Conference on Learning Representations*, 2022.
- Hongye Cao, Fan Feng, Meng Fang, Shaokang Dong, Jing Huo, and Yang Gao. Towards empowerment gain through causal structure learning in model-based rl. In *ICML 2024 Workshop: Foundations of Reinforcement Learning and Control—Connections and Perspectives*, 2024.
- Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Variational empowerment as representation learning for goal-based reinforcement learning. *arXiv preprint arXiv:2106.01404*, 2021.
- ZH Deng, J Jiang, G Long, and C Zhang. Causal reinforcement learning: A survey. *Transactions on Machine Learning Research*, 2023.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2018.
- Fan Feng and Sara Magliacane. Learning dynamic attribute-factored world models for efficient multi-object reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2023.
- Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. *Advances in neural information processing systems*, 14, 2001.
- Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms for factored mdps. *Journal of Artificial Intelligence Research*, 19:399–468, 2003.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- Biwei Huang, Fan Feng, Chaochao Lu, Sara Magliacane, and Kun Zhang. Adarl: What, where, and how to adapt in transfer reinforcement learning. In *International Conference on Learning Representations*, 2022a.
- Biwei Huang, Chaochao Lu, Liu Leqi, José Miguel Hernández-Lobato, Clark Glymour, Bernhard Schölkopf, and Kun Zhang. Action-sufficient state representation learning for control with structural constraints. In *International Conference on Machine Learning*, pp. 9260–9279. PMLR, 2022b.
- Tianying Ji, Yongyuan Liang, Yan Zeng, Yu Luo, Guowei Xu, Jiawei Guo, Ruijie Zheng, Furong Huang, Fuchun Sun, and Huazhe Xu. Ace: Off-policy actor-critic with causality-aware entropy regularization. In *Forty-first International Conference on Machine Learning*, 2024a.
- Tianying Ji, Yu Luo, Fuchun Sun, Xianyuan Zhan, Jianwei Zhang, and Huazhe Xu. Seizing serendipity: Exploiting the value of past success in off-policy actor-critic. In *Forty-first International Conference on Machine Learning*, 2024b.
- Tobias Jung, Daniel Polani, and Peter Stone. Empowerment for continuous agent—environment systems. *Adaptive Behavior*, 19(1):16–39, 2011.
- Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *IJCAI*, volume 16, pp. 740–747, 1999.
- Alexander S Klyubin, Daniel Polani, and Chrystopher L Nehaniv. Empowerment: A universal agent-centric measure of control. In *2005 ieee congress on evolutionary computation*, volume 1, pp. 128–135. IEEE, 2005.

- Oliver Kroemer, Scott Niekum, and George Konidaris. A review of robot learning for manipulation: Challenges, representations, and algorithms. *Journal of machine learning research*, 22(30):1–82, 2021.
- Felix Leibfried, Sergio Pascual-Diaz, and Jordi Grau-Moya. A unified bellman optimality principle combining reward maximization and empowerment. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yuren Liu, Biwei Huang, Zhengmao Zhu, Honglong Tian, Mingming Gong, Yang Yu, and Kun Zhang. Learning world models with identifiable factorization. *Advances in Neural Information Processing Systems*, 36, 2024.
- Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsically motivated reinforcement learning. *Advances in neural information processing systems*, 28, 2015.
- Mirco Mutti, Riccardo De Santi, Emanuele Rossi, Juan Felipe Calderon, Michael Bronstein, and Marcello Restelli. Provably efficient causal model-based reinforcement learning for systematic generalization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9251–9259, 2023.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Silviu Pitis, Elliot Creager, and Animesh Garg. Counterfactual data augmentation using locally factored dynamics. *Advances in Neural Information Processing Systems*, 33:3976–3990, 2020.
- Silviu Pitis, Elliot Creager, Ajay Mandlekar, and Animesh Garg. Mocoda: Model-based counterfactual data augmentation. *Advances in Neural Information Processing Systems*, 35:18143–18156, 2022.
- Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *Robotics: Science and Systems XIV*, 2018.
- Jonathan Richens and Tom Everitt. Robust agents learn causal world models. *arXiv preprint arXiv:2402.10877*, 2024.
- Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment—an introduction. *Guided Self-Organization: Inception*, pp. 67–114, 2014.
- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9339–9347, 2019.
- Maximilian Seitzer, Bernhard Schölkopf, and Georg Martius. Causal influence detection for improving efficiency in reinforcement learning. *Advances in Neural Information Processing Systems*, 34: 22905–22918, 2021.
- Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara, Takashi Washio, Patrik O Hoyer, Kenneth Bollen, and Patrik Hoyer. Directlingam: A direct method for learning a linear non-gaussian structural equation model. *Journal of Machine Learning Research-JMLR*, 12(Apr):1225–1248, 2011.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press, 2001.
- Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033. IEEE, 2012.
- Núria Armengol Urpí, Marco Bagatella, Marin Vlastelica, and Georg Martius. Causal action influence aware counterfactual data augmentation. In *Forty-first International Conference on Machine Learning*, 2024.
- Zizhao Wang, Xuesu Xiao, Yuke Zhu, and Peter Stone. Task-independent causal state abstraction. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, Robot Learning workshop*, 2021.
- Zizhao Wang, Xuesu Xiao, Zifan Xu, Yuke Zhu, and Peter Stone. Causal dynamics learning for task-independent state abstraction. In *International Conference on Machine Learning*, pp. 23151–23180. PMLR, 2022.
- Zizhao Wang, Caroline Wang, Xuesu Xiao, Yuke Zhu, and Peter Stone. Building minimal and reusable causal state abstractions for reinforcement learning. *arXiv preprint arXiv:2401.12497*, 2024.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
- Yanjie Ze, Yuyao Liu, Ruizhe Shi, Jiaxin Qin, Zhecheng Yuan, Jiashun Wang, and Huazhe Xu. H-index: Visual reinforcement learning with hand-informed representations for dexterous manipulation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yan Zeng, Ruichu Cai, Fuchun Sun, Libo Huang, and Zhifeng Hao. A survey on causal reinforcement learning. *arXiv preprint arXiv:2302.05209*, 2023.
- Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. Interpretable reward redistribution in reinforcement learning: a causal approach. *Advances in Neural Information Processing Systems*, 36, 2024.

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Work</b>	<b>3</b>
2.1	Causal RL . . . . .	3
2.2	Empowerment in RL . . . . .	3
<b>3</b>	<b>Preliminaries</b>	<b>3</b>
3.1	Markov Decision Process . . . . .	3
3.2	Structural Causal Model . . . . .	4
3.3	Empowerment in RL . . . . .	4
<b>4</b>	<b>CIP</b>	<b>4</b>
4.1	Counterfactual Data Augmentation . . . . .	5
4.2	Causal Action Prioritization Through Empowerment . . . . .	5
<b>5</b>	<b>Experiments</b>	<b>7</b>
5.1	Experimental setup . . . . .	7
5.2	Main Results . . . . .	8
5.3	Property Analysis . . . . .	9
<b>6</b>	<b>Conclusion</b>	<b>10</b>
<b>A</b>	<b>Broader Impact</b>	<b>16</b>
<b>B</b>	<b>Assumptions and Propositions</b>	<b>16</b>
<b>C</b>	<b>Extensive Related Work</b>	<b>17</b>
<b>D</b>	<b>Details on Experimental Design and Results</b>	<b>18</b>
D.1	Experimental setup . . . . .	18
D.2	Full Results . . . . .	18
D.2.1	Effectiveness in robot arm manipulation . . . . .	18
D.2.2	Effectiveness in spare reward settings . . . . .	18
D.2.3	Effectiveness in locomotion . . . . .	20
D.2.4	Effectiveness in pixel-based tasks . . . . .	23
D.3	Property Analysis . . . . .	23
D.3.1	Analysis for replacing counterfactual data augmentation . . . . .	23
D.3.2	Extensive ablation study . . . . .	24
D.3.3	Hyperparameter analysis . . . . .	25
<b>E</b>	<b>Details on the Proposed Framework</b>	<b>28</b>



756	<b>F Experimental Platforms and Licenses</b>	<b>28</b>
757		
758	F.1 Experimental platforms . . . . .	28
759	F.2 Licenses . . . . .	28
760		
761		
762		
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

## A BROADER IMPACT

To avoid blind exploration and improve sample efficiency, we propose **CIP** for efficient reinforcement learning. **CIP** leverages the causal relationships among states, actions, and rewards to prioritize causal information for efficient policy learning. **CIP** first learns a causal matrix between states and rewards to execute counterfactual data augmentation, prioritizing important state features without additional environmental interactions. Subsequently, it learns a causal reweight matrix between actions and rewards to prioritize causally-informed behaviors. We then introduce a causal action empowerment term into the learning objective to enhance the controllability. By prioritizing the causal information, **CIP** enables agents to focus on behaviors that have causally significant effects on their tasks. **CIP** offers substantial broader impact by prioritizing causal information through individual assessment of how different factors contribute to rewards. Our novel empowerment learning objective achieves efficient policy optimization by leveraging entropy via the policy and learned inverse dynamics model. This approach shows promise for extension into research frameworks centered on maximum entropy algorithms.

Despite its strengths, **CIP** has limitations beyond its reliance on the method DirectLiNGAM. There’s potential to explore alternative causal discovery techniques for more robust relationship mapping. Moreover, analyzing inter-entity causal connections could lead to better disentanglement of diverse behaviors. Our future work will investigate a range of causal discovery methods to refine our approach. We aim to extend **CIP** to model-based RL frameworks, focusing on building causal world models to enhance generalization.

## B ASSUMPTIONS AND PROPOSITIONS

**Assumption 1** (*d-separation* (Pearl, 2009)) *d-separation is a graphical criterion used to determine, from a given causal graph, if a set of variables  $X$  is conditionally independent of another set  $Y$ , given a third set of variables  $Z$ . In a directed acyclic graph (DAG)  $\mathcal{G}$ , a path between nodes  $n_1$  and  $n_m$  is said to be blocked by a set  $S$  if there exists a node  $n_k$ , for  $k = 2, \dots, m-1$ , that satisfies one of the following two conditions:*

(i)  $n_k \in S$ , and the path between  $n_{k-1}$  and  $n_{k+1}$  forms  $(n_{k-1} \rightarrow n_k \rightarrow n_{k+1})$ ,  $(n_{k-1} \leftarrow n_k \leftarrow n_{k+1})$ , or  $(n_{k-1} \leftarrow n_k \rightarrow n_{k+1})$ .

(ii) Neither  $n_k$  nor any of its descendants is in  $S$ , and the path between  $n_{k-1}$  and  $n_{k+1}$  forms  $(n_{k-1} \rightarrow n_k \leftarrow n_{k+1})$ .

In a DAG, we say that two nodes  $n_a$  and  $n_b$  are *d-separated* by a third node  $n_c$  if every path between nodes  $n_a$  and  $n_b$  is blocked by  $n_c$ , denoted as  $n_a \perp\!\!\!\perp n_b | n_c$ .

**Assumption 2** (*Global Markov Condition* (Spirtes et al., 2001; Pearl, 2009)) *The state is fully observable and the dynamics is Markovian. The distribution  $p$  over a set of variables  $\mathcal{V} = (s_t^1, \dots, s_t^d, a_t^1, \dots, a_t^d, r_t)^T$  satisfies the global Markov condition on the graph if for any partition  $(\mathcal{S}, \mathcal{A}, \mathcal{R})$  in  $\mathcal{V}$  such that if  $\mathcal{A}$  d-separates  $\mathcal{S}$  from  $\mathcal{R}$ , then  $p(\mathcal{S}, \mathcal{R} | \mathcal{A}) = p(\mathcal{S} | \mathcal{A}) \cdot p(\mathcal{R} | \mathcal{A})$*

**Assumption 3** (*Faithfulness Assumption* (Spirtes et al., 2001; Pearl, 2009)) *For a set of variables  $\mathcal{V} = (s_t^1, \dots, s_t^d, a_t^1, \dots, a_t^d, r_t)^T$ , there are no independencies between variables that are not implied by the Markovian Condition.*

**Assumption 4** *Under the assumptions that the causal graph is Markov and faithful to the observations, the edge  $s_t^i \rightarrow s_{t+1}^i$  exists for all state variables  $s^i$ .*

**Assumption 5** *No simultaneous or backward edges in time.*

**Proposition 1** *Under the assumptions that the causal graph is Markov and faithful to the observations, there exists an edge from  $a_t^i \rightarrow r_t$  if and only if  $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ .*

*Proof:* We proceed by proving both directions of the if and only if statement.

( $\Rightarrow$ ) Suppose there exists an edge from  $a_t^i$  to  $r_t$ . We prove that  $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$  by contradiction. Assume  $a_t^i \perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ . By the faithfulness assumption, this independence must be reflected in

Table 2: Categorization of different causal RL methods with two different causal relationship of state-to-reward (state-reward) and action-to-reward (action-reward).

Problem domain	Task type	Method	Causal relationship	
			state-reward	action-reward
<b>Single-task</b>	manipulation; locomotion	ACE (Ji et al., 2024a)	✗	✓
	manipulation; locomotion	IFactor (Liu et al., 2024)	✓	✗
	manipulation	CAI (Seitzer et al., 2021)	✗	✗
<b>Generalization</b>	manipulation	CDL (Wang et al., 2022)	✗	✗
	locomotion	AdaRL (Huang et al., 2022a)	✓	✓
	manipulation; locomotion	CBM (Wang et al., 2024)	✓	✗
<b>Augmentation</b>	manipulation	CAIAC (Urpí et al., 2024)	✗	✗
	manipulation	CoDA (Pitis et al., 2020)	✗	✗
	manipulation	MoCoDA (Pitis et al., 2022)	✗	✗

the graph structure. However, this implies the absence of a directed path from  $a_t^i$  to  $r_t$ , contradicting the existence of the edge. Thus,  $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ .

( $\Leftarrow$ ) Now, suppose  $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ . We prove the existence of an edge from  $a_t^i$  to  $r_t$  by contradiction. Assume no such edge exists. By the Markov assumption, the absence of this edge implies  $a_t^i \perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ , contradicting our initial supposition. Therefore, an edge from  $a_t^i$  to  $r_t$  must exist. Thus, we have shown that an edge from  $a_t^i$  to  $r_t$  exists if and only if  $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ , completing the proof.

**Proposition 2** *Under the assumptions that the causal graph is Markov and faithful to the observations, there exists an edge from  $s_t^i \rightarrow r_t$  if and only if  $s_t^i \not\perp\!\!\!\perp r_t | \{a_t, s_t \setminus r_t\}$ .*

The proof of Proposition 2 follows a similar line of reasoning as that of Proposition 1.

**Theorem 1** *Based on above 5 assumptions and 2 propositions, suppose  $s_t, a_t, s_t$  follow the factored MDP reward function Eq. 2, the causal matrices  $M^{s \rightarrow r}$  and  $M^{a \rightarrow r}$  are identifiable.*

## C EXTENSIVE RELATED WORK

We categorize existing causal RL approaches based on problem domains and task types, providing a systematic analysis of how different methods explore causal relationships between states, actions, and rewards, as illustrated in Table 2.

In the single-task learning domain, methods such as ACE (Ji et al., 2024a) and IFactor (Liu et al., 2024) have shown success in learning policies for manipulation and locomotion tasks. However, both approaches are limited by focusing on a single reward-guided causal relationship. Regarding generalization, AdaRL (Huang et al., 2022a) effectively leverages both state-reward and action-reward causal relationships. However, AdaRL focuses primarily on applying causal inference to address generalization challenges in locomotion tasks. Its application is limited to locomotion tasks, leaving more complex manipulation tasks unaddressed. Since our work focuses on the single-task problem domain, we do not provide a direct comparison with AdaRL. Conversely, CBM (Wang et al., 2024) considers the causal relationship between states and rewards but overlooks the causal link between actions and rewards. In the problem domain of counterfactual data augmentation, current causal RL methods (Urpí et al., 2024; Pitis et al., 2020; 2022) have not yet explored the inference and utilization of both causal relationships.

In summary, current research on reward-guided causal discovery remains incomplete and lacks validation across a broader spectrum of tasks. This gap underscores the need for more comprehensive investigation and application in the field of causal reinforcement learning.

## D DETAILS ON EXPERIMENTAL DESIGN AND RESULTS

### D.1 EXPERIMENTAL SETUP

We present the detailed hyperparameter settings of the proposed method **CIP** across all 5 environments in Table 3. Additionally, the Q-value and V-value networks are used MLP with 512 hidden size. And the policy network is the Gaussian MLP with 512 hidden size. Moreover, we set the target update interval of 2. For fair comparison, the hyperparameters of the baseline methods (SAC (Haarnoja et al., 2018), BAC (Ji et al., 2024b), ACE (Ji et al., 2024a)) follow the same settings in the experiments.

For pixel-based DMControl environments, we employ IFactor (Liu et al., 2024) to encode latent states and integrate the **CIP** framework for policy learning. We utilize the  $s_t^{\bar{}}$  state features in IFactor as uncontrollable states unrelated to rewards to execute counterfactual data augmentation. Furthermore, for simplicity, we maximize the mutual information between future states and actions to facilitate empowerment. All parameter settings in these three tasks adhere to those specified in IFactor. Additionally, We use the same background video for the comparison.

Table 3: Hyperparameter settings of **CIP** in 5 environments

Hyperparameter	Environment				
	Meta-World	Sparse	MuJoCo	DMControl	Adroit Hand
batch size	512	512	256	512	256
hidden size	1024	1024	256	1024	256
Q-value network hidden size			512		
V-value network hidden size			512		
policy network hidden size			512		
learning step			1000000		
replay size			1000000		
causal sample size			10000		
gamma			0.99		
learning rate			0.0003		
update interval			2		

### D.2 FULL RESULTS

#### D.2.1 EFFECTIVENESS IN ROBOT ARM MANIPULATION

Figure 10 presents the learning curves for all 17 manipulation skill tasks within the Meta-World environment. The **CIP** framework demonstrates superior learning outcomes and efficiency compared to the three baseline methods, despite exhibiting minor instabilities in the basketball and dial-turn tasks. Notably, **CIP** achieves a 100% success rate in more complex tasks, such as pick-place-wall and assembly. The visualization results presented in Figures 11 and 12 further demonstrate **CIP**’s ability to effectively and efficiently complete tasks, even in high-dimensional action spaces such as the Adroit Hand environment.

In the hammer task, **CIP** allows the robot arm to execute reach and pick actions with precision, enabling it to accurately identify the nail’s position and successfully perform the hammering action. In the Adroit Hand door task, **CIP** effectively controls the complex joints to grasp the doorknob and applies the appropriate force to twist it, thereby opening the door.

These findings affirm the effectiveness of **CIP** in robot arm manipulation skill learning, highlighting its capacity to enhance sample efficiency while mitigating the risks associated with blind exploration.

#### D.2.2 EFFECTIVENESS IN SPARE REWARD SETTINGS

Figure 13 presents the learning curves for all three sparse reward setting tasks within the Meta-World environment, while Figure 14 showcases their corresponding visualization trajectories. These findings

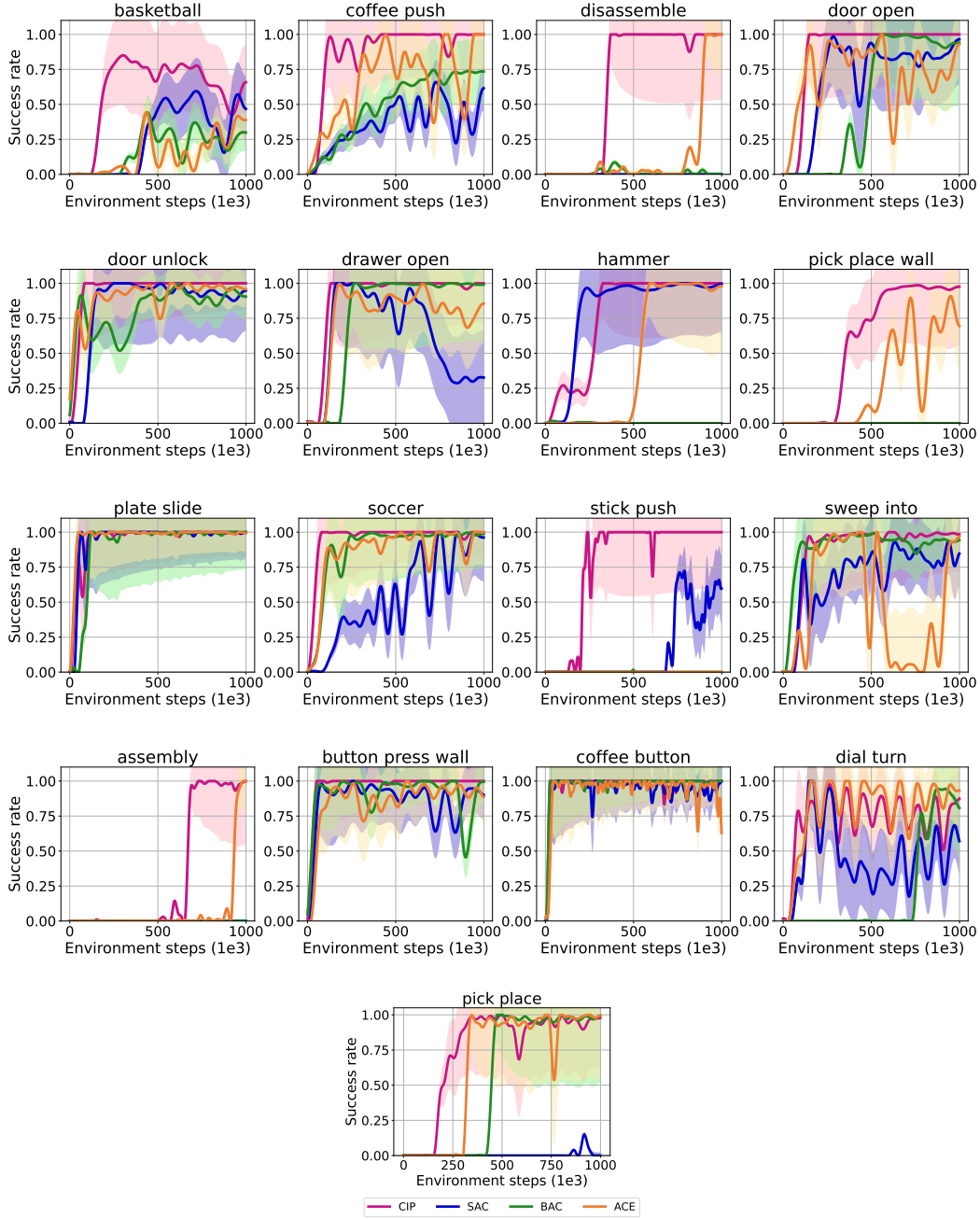


Figure 10: Experimental results across 17 manipulation skill learning tasks in Meta-World.

reveal that **CIP** not only achieves superior learning efficiency but also adeptly executes critical actions necessary for task completion, such as opening the door and window and maneuvering the node to the target place.

These results substantiate the effectiveness of **CIP** in sparse reward scenarios. The counterfactual data augmentation process prioritizes salient state information, effectively filtering out irrelevant factors that could hinder learning. Meanwhile, causal action empowerment enhances policy controllability by focusing on actions that are causally linked to desired outcomes. This dual approach not only accelerates the learning process but also fosters a more robust policy capable of navigating the complexities inherent in sparse reward settings. Overall, these findings underscore **CIP**'s potential to significantly improve performance in challenging environments characterized by limited feedback.



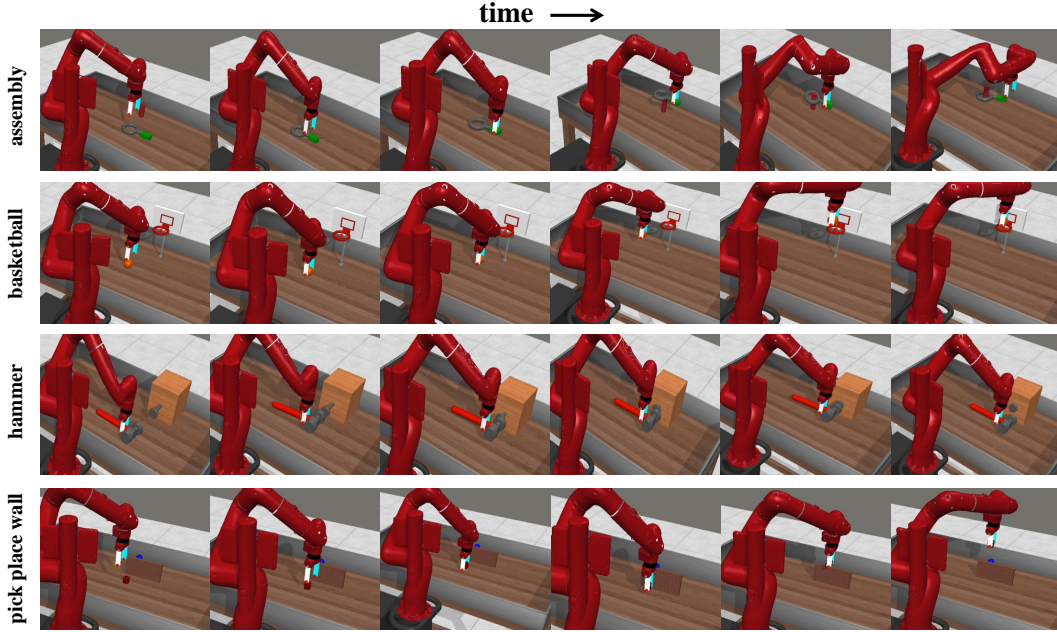


Figure 11: Visualization trajectories of 4 manipulation skill learning tasks in Meta-World environment.

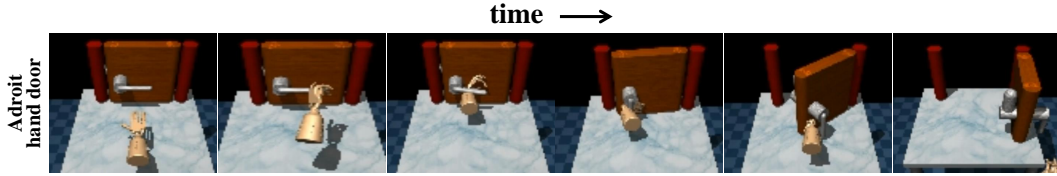


Figure 12: Visualization trajectory of Adroit Hand door open task.

### D.2.3 EFFECTIVENESS IN LOCOMOTION

We further evaluate **CIP** in 15 locomotion tasks in DMControl and MuJoCo environments. Figure 15 presents the learning curves, while Figure 16 showcases the corresponding visualization trajectories in 4 specific tasks. A comprehensive analysis indicates that **CIP** achieves faster learning efficiency and greater stability compared to ACE and SAC, while demonstrating comparable policy learning performance to BAC, which is known for its proficiency in control tasks. The visualization results reveal that **CIP** effectively executes running and walking actions in complex humanoid scenarios.

These findings collectively underscore the efficacy of **CIP** in locomotion tasks, highlighting its potential to advance the state-of-the-art in reinforcement learning for intricate motor control problems. The method’s success across varied environments suggests a robust framework that could generalize effectively to other challenging domains within robotics and control systems, paving the way for future research and applications in these areas.

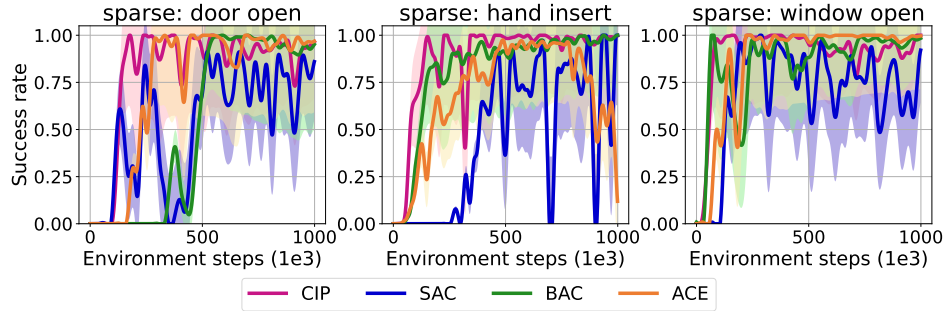


Figure 13: Experimental results across 3 manipulation skill learning tasks in sparse reward settings of Meta-World environment.

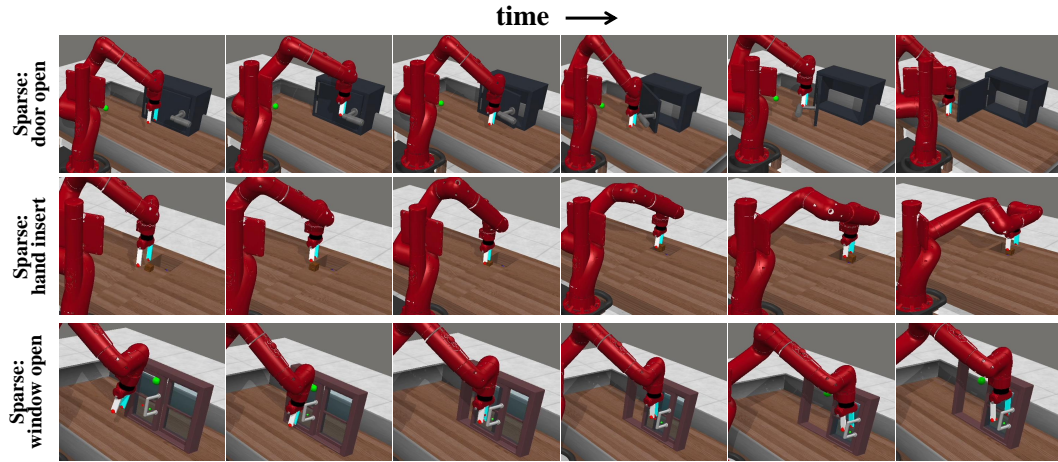


Figure 14: Visualization trajectories of 3 manipulation skill learning tasks in sparse reward settings of Meta-World environment.

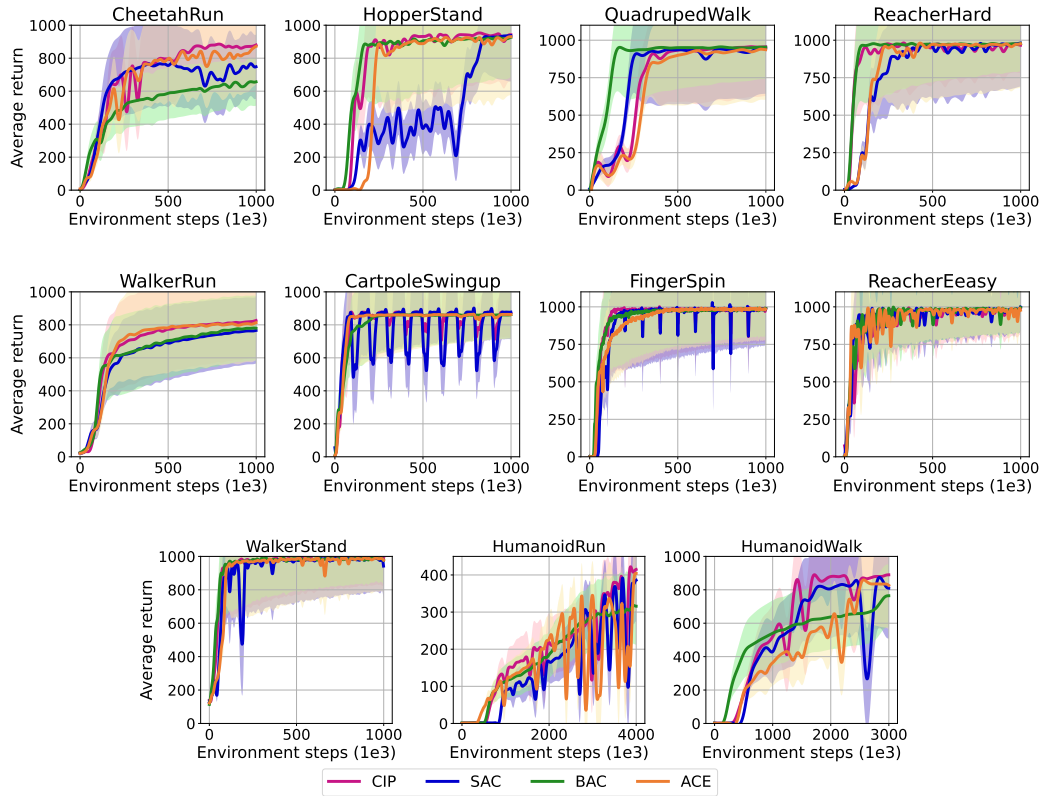


Figure 15: Experimental results across 11 locomotion tasks in DMControl environment.

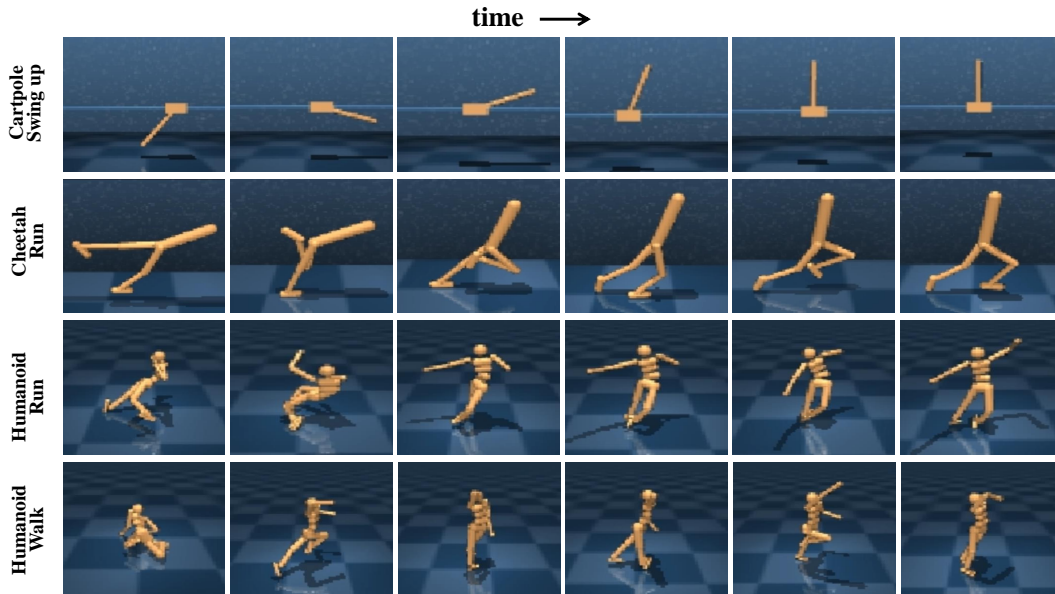


Figure 16: Visualization trajectories of 4 locomotion tasks in DMControl environment.

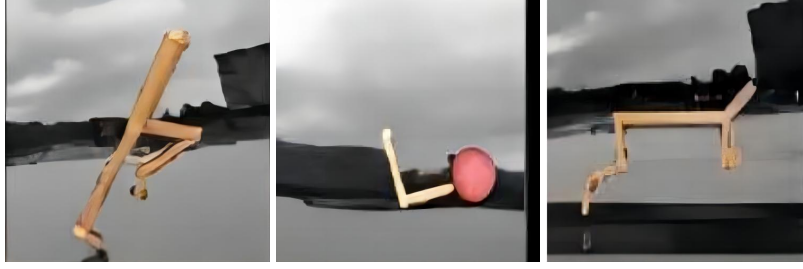


Figure 17: The DMControl environment of 3 pixel-based tasks (Walker Walk, Cheetah Run, Reacher Easy).

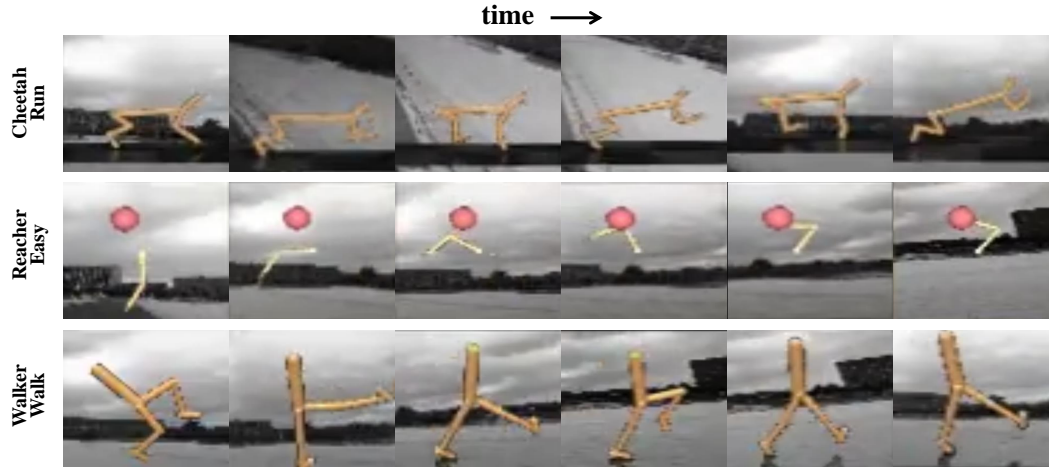


Figure 18: Visualization trajectories in 3 pixel-based locomotion tasks of DMControl environment with video backgrounds as distractors.

#### D.2.4 EFFECTIVENESS IN PIXEL-BASED TASKS

To further validate the effectiveness of our proposed framework in pixel-based environments, we evaluated **CIP** on three DMControl pixel-based tasks. We leverage IFactor for latent state processing and differentiation of uncontrollable state features to execute counterfactual data augmentation, alongside maximizing the mutual information between future states and actions for empowerment.

Figure 6 presents the learning curves, while Figure 18 shows the visualization trajectories. The proposed framework exhibits enhanced policy learning performance and effectively mitigates interference from background video, facilitating efficient locomotion. These findings reinforce the effectiveness and extensibility of our causal information prioritization framework, highlighting its potential to improve learning in complex, pixel-based environments.

### D.3 PROPERTY ANALYSIS

#### D.3.1 ANALYSIS FOR REPLACING COUNTERFACTUAL DATA AUGMENTATION

In **CIP**, we exploit the causal relationship between states and rewards to perform counterfactual data augmentation on irrelevant state features, thus prioritizing critical state information. We compare this approach with an alternative method: masking irrelevant state features to achieve state abstraction for subsequent causal action empowerment and policy learning. To evaluate the efficacy of both approaches, we conduct experiments with **CIP** with counterfactual data augmentation (**CIP** w/i Cda) and **CIP** with causally-informed states (**CIP** w/i Cs) across three distinct environments.

Figure 19 illustrates comparative results for four manipulation skill learning tasks in the Meta-World environment. Both **CIP** variants achieve 100% task success rates with high sample efficiency, vali-



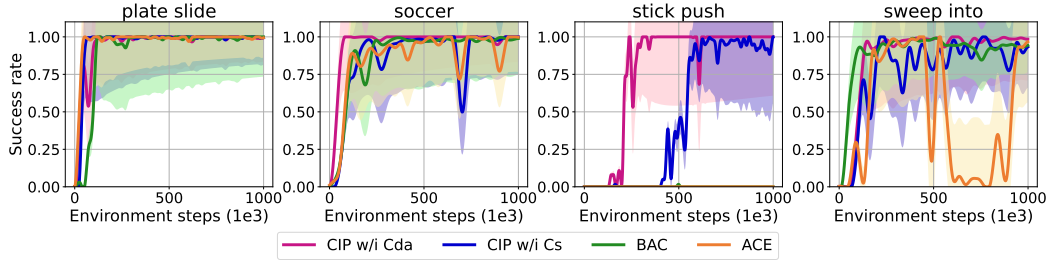


Figure 19: Experimental results in 4 manipulation skill learning tasks of Meta-World environment. w/i stands for with.

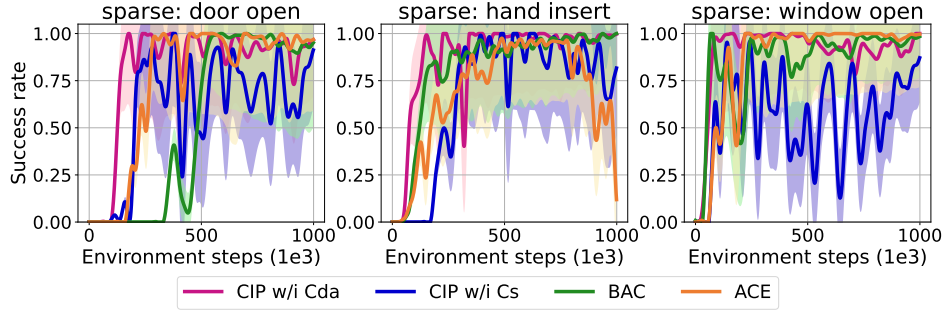


Figure 20: Experimental results in 3 manipulation skill learning tasks of Meta-World environment with sparse reward settings.

dating their effectiveness. Notably, **CIP w/i Cda** exhibits superior learning efficiency compared to **CIP w/i Cs**, underscoring the value of our counterfactual data augmentation approach in enhancing training data without additional environmental interactions. In three sparse reward setting tasks (Figure 20), **CIP w/i Cda** demonstrates superior policy performance. Further experiments across four locomotion environment tasks corroborate these findings, consistently favoring the counterfactual data augmentation approach. These comprehensive experimental results strongly support the effectiveness and significance of incorporating counterfactual data augmentation in **CIP**, highlighting its potential to enhance reinforcement learning across diverse task domains.

### D.3.2 EXTENSIVE ABLATION STUDY

**Robot arm manipulation** The ablation study results in the Meta-World and Adroit Hand environments are presented in Figure 22. The findings indicate that **CIP** without counterfactual data augmentation exhibits reduced learning efficiency and is unable to successfully complete tasks such as pick-and-place. This underscores the importance of incorporating counterfactual data augmentation, which prioritizes causal state information, to enhance learning efficiency by mitigating the influence of irrelevant state information and preventing policy divergence.

Furthermore, **CIP** without causal action empowerment demonstrates a significant decline in policy performance across robot arm manipulation tasks. In complex scenarios, such as Adroit Hand door opening and assembly, it fails to learn effective strategies for task completion. This outcome further corroborates the efficacy of the proposed causal action empowerment mechanism, as prioritizing causally informed actions facilitates more efficient exploration of the environment, ultimately enabling successful policy learning.

**Sparse reward settings** Figure 22 presents the results of the ablation study conducted across three sparse reward setting tasks. These findings underscore the substantial influence of causal action empowerment on the efficacy of policy learning, demonstrating its critical role in enhancing performance in challenging environments. Additionally, the incorporation of counterfactual data augmentation proves effective in mitigating the need for additional environmental interactions, thereby significantly improving sample efficiency. This approach not only facilitates more rapid learning but



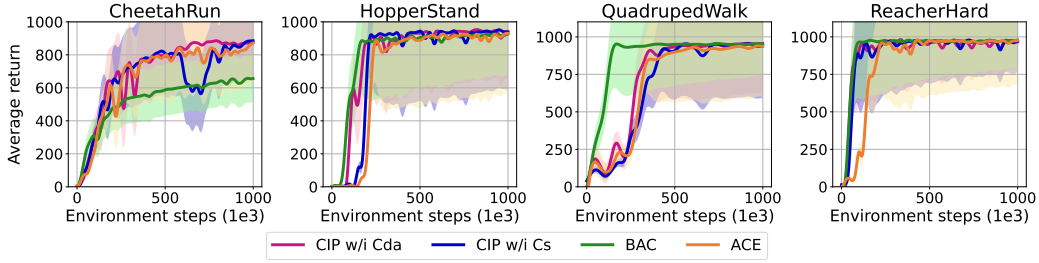


Figure 21: Experimental results in 4 locomotion tasks of DMControl environment.

also ensures that the agent can effectively navigate sparse reward scenarios by focusing on the most relevant causal information.

**Locomotion** We further conducted ablation experiments on locomotion tasks. The experimental results in the MuJoCo environment are shown in Figure 23, where it is evident that the performance of **CIP** without causal action empowerment declines significantly. Similarly, **CIP** without counterfactual data augmentation also exhibits reduced learning efficiency. Notably, in the 11 DMControl tasks, the decline in performance for **CIP** without causal action empowerment is particularly pronounced.

These experimental results further validate the effectiveness of our proposed method, which systematically analyzes the causal relationships between states, actions, and rewards. This analysis enables the execution of counterfactual data augmentation to avoid interference from irrelevant factors while prioritizing important state information. Subsequently, by leveraging the causal relationships between actions and rewards, we reweight actions to prioritize causally informed actions, thereby enhancing the agent’s controllability and overall learning efficacy.

### D.3.3 HYPERPARAMETER ANALYSIS

We conduct a detailed analysis of the hyperparameters associated with the causal update interval ( $I$ ) and sample size within the **CIP** framework. The experimental results for four distinct tasks are illustrated in Figure 25. Across all tasks, **CIP** demonstrates optimal performance with a causal update interval of  $I = 2$  and a sample size of 10,000.

Our findings suggest that while a reduction in the causal update interval can lead to improved performance, it may also result in heightened computational costs. Additionally, we observe that higher update frequencies and increased sample sizes introduce greater instability, which significantly raises computational demands. This analysis underscores the importance of carefully balancing hyperparameter settings to optimize both performance and efficiency within the **CIP**.

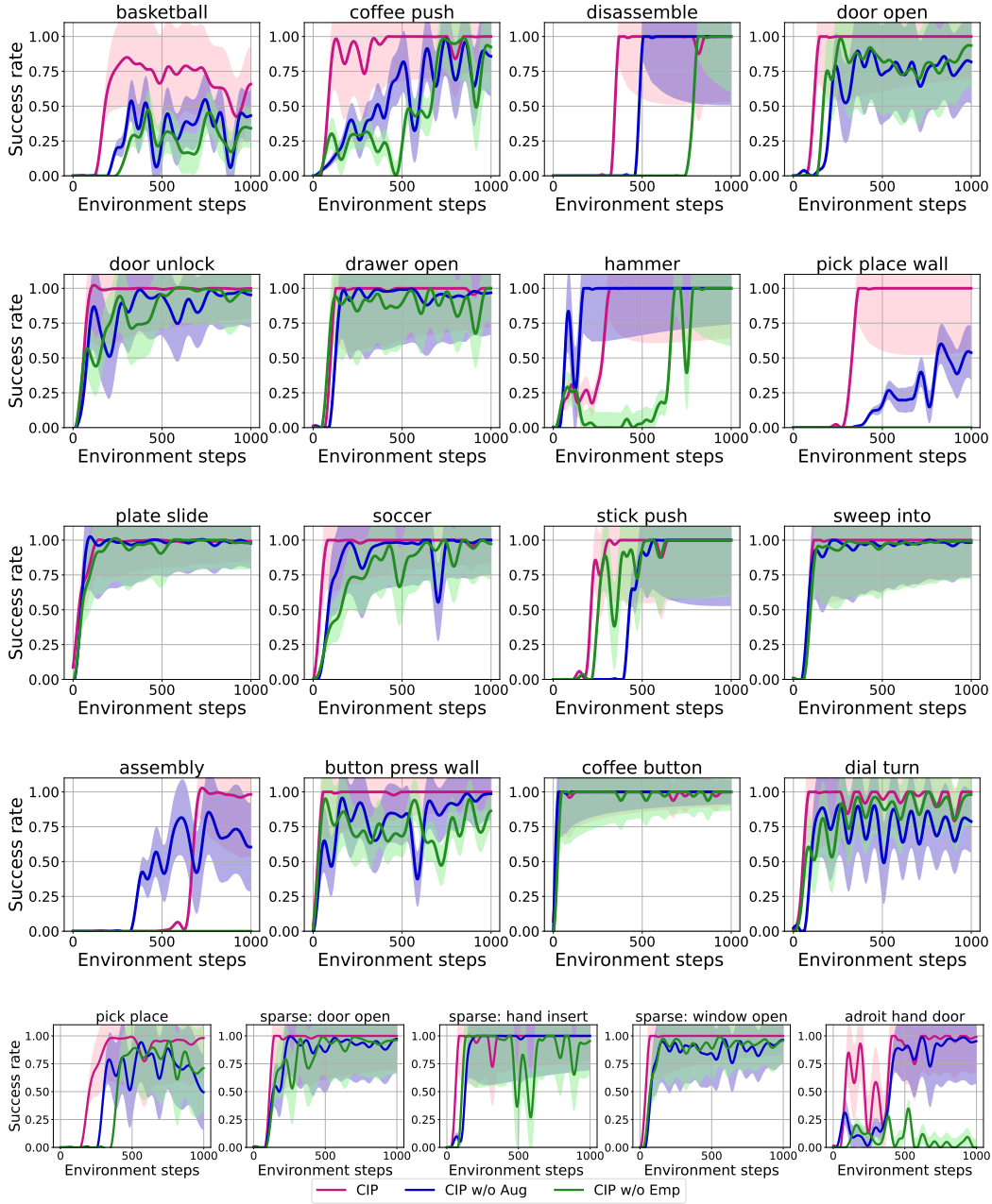


Figure 22: Ablation results across 21 manipulation skill learning tasks in Meta-World including sparse reward settings and adroit hand.

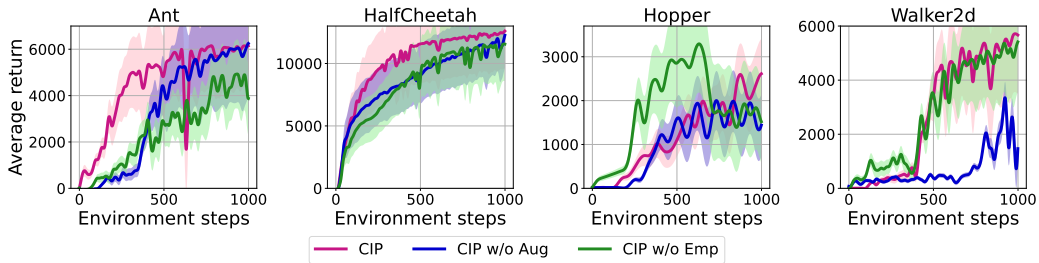


Figure 23: Ablation results across 4 locomotion tasks in MuJoCo environment.

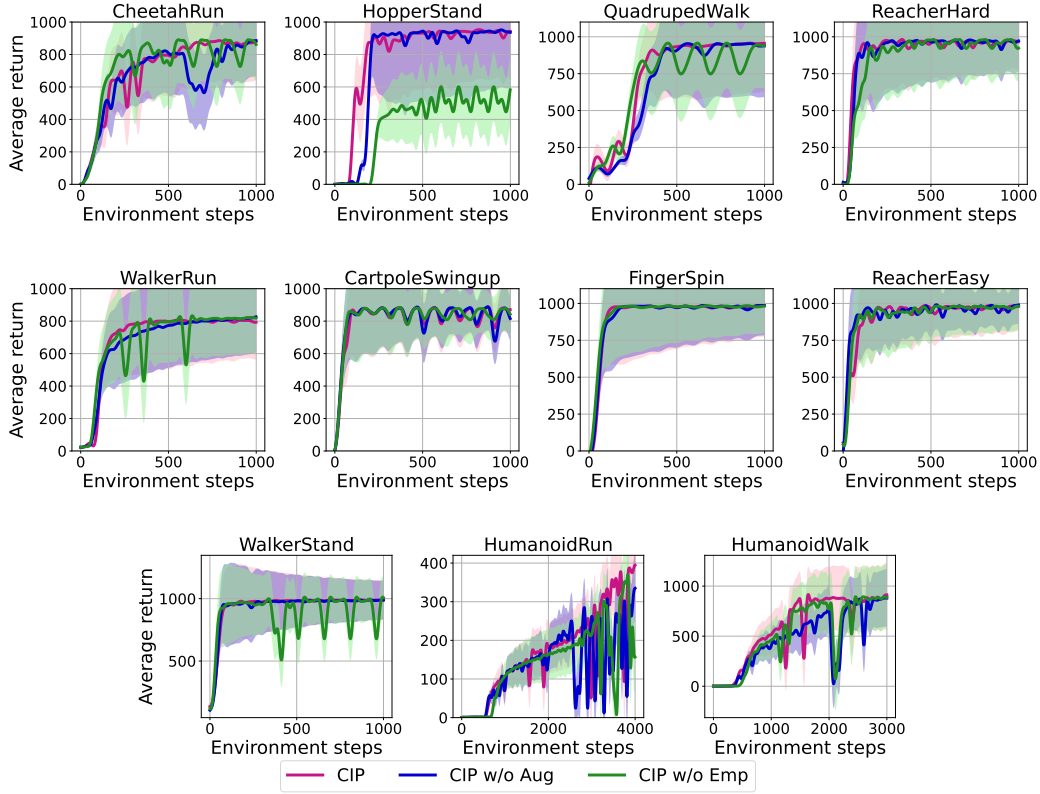


Figure 24: Ablation results across 11 locomotion tasks in DMControl environment.

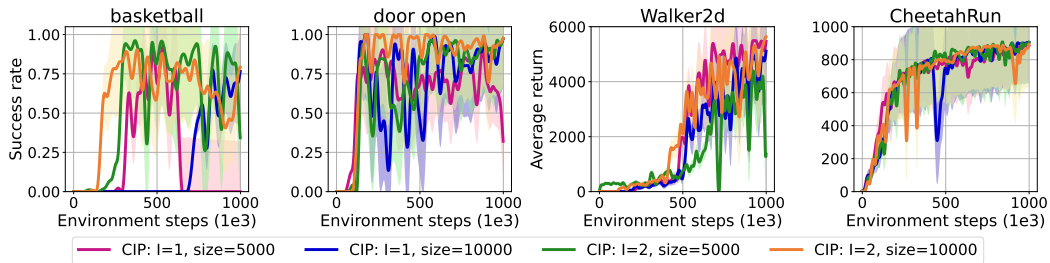


Figure 25: Hyperparameter study. Learning curves of **CIP** with different hyperparameter settings. The shaded regions are the standard deviation of each policy.

## E DETAILS ON THE PROPOSED FRAMEWORK

Algorithm 1 lists the full pipeline of **CIP** below.

---

### Algorithm 1 Causal information prioritization for efficient RL

---

**Input:**  $Q$  network  $Q_{\pi_c}$ , policy network  $\pi_c$ , inverse dynamics model  $\phi_c$  with  $Q$  network  $Q_{\phi_c}$ , replay buffer  $\mathcal{D}$ , local causal buffer  $\mathcal{D}_c$ , causal update interval  $I$ , causal matrix  $M^{a \rightarrow s}$  and  $M^{a \rightarrow r}$ .

**for** each environment step  $t$  **do**  
 Collect data with  $\pi_\theta$  from real environment  
 Add to replay buffer  $\mathcal{D}$  and local buffer  $\mathcal{D}_c$   
**end for**

#### Step 1: Counterfactual data augmentation

**if** every  $I$  environment step **then**  
 Sample transitions  $\mathcal{D}_s$  from local buffer  $\mathcal{D}_c$   
 Learn causal mask matrix  $M^{a \rightarrow r}$  with  $\{(s, a, r, s')\}^{|\mathcal{D}_s|}$  for causal state prioritization  
 Compute uncontrollable set  $\mathcal{U}_s$  followed by Eq. 4  
 Sample  $(s, a, r, s') \in \mathcal{D}_s$   
**for**  $s^i \in \mathcal{U}_s$  **do**  
 Sample  $(\hat{s}, \hat{a}, \hat{r}, \hat{s}') \sim \mathcal{D}_s$   
**if** state  $\hat{s}^i \in \mathcal{U}_s$  **then**  
 Construct a counterfactual transition  $(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}')$  by swapping  $(s^i, s'^i)$  with  $(\hat{s}^i, \hat{s}'^i)$   
 Add  $(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}')$  to local buffer  $\mathcal{D}_c$   
**end if**  
**end for**  
**end if**

#### Step 2: Causal weighted matrix learning

**if** every  $I$  environment step **then**  
 Sample transitions  $\mathcal{D}_a$  from local buffer  $\mathcal{D}_c$   
 Learn causal weighted matrix  $M^{a \rightarrow r}$  with  $\{(s, a, r, s')\}^{|\mathcal{D}_a|}$  for causal action prioritization  
**end if**

#### Step 3: Policy optimization with causal action empowerment

**for** each gradient step **do**  
 Sample  $N$  transitions  $(s, a, r, s')$  from  $\mathcal{D}$   
 Compute causal action empowerment followed by Eq. 8.  
 Calculate the target  $Q_{\phi_c}$  value  
 Update  $Q_{\phi_c}$  by  $\min_{\phi_c} (\mathcal{T}Q_{\phi_c} - Q_{\phi_c})^2$   
 Update  $\phi_c$  by  $\max(Q_{\phi_c}(s, a))$   
 Calculate the target  $Q_{\pi_c}$  value  
 Update  $Q_{\pi_c}$  by  $\min_{\pi_c} (\mathcal{T}_c Q_{\pi_c} - Q_{\pi_c})^2$   
 Update  $\pi_c$  by  $\max_c (Q_{\pi_c}(s, a) + \mathcal{E}_{\pi_c}(s))$   
**end for**

---

## F EXPERIMENTAL PLATFORMS AND LICENSES

### F.1 EXPERIMENTAL PLATFORMS

All experiments of this approach are implemented on 2 Intel(R) Xeon(R) Gold 6430 and 2 NVIDIA Tesla A800 GPUs.

### F.2 LICENSES

In our code, we have utilized the following libraries, each covered by its respective license agreements:

- PyTorch (BSD 3-Clause "New" or "Revised" License)
- Numpy (BSD 3-Clause "New" or "Revised" License)
- Tensorflow (Apache License 2.0)
- Meta-World (MIT License)
- MuJoCo (Apache License 2.0)
- Deep Mind Control (Apache License 2.0)
- Adroit Hand (Creative Commons License 3.0)